

# Neural Acceleration for GPU Throughput Processors

Amir Yazdanbakhsh      Jongse Park      Hardik Sharma  
Pejman Lotfi-Kamran<sup>†</sup>      Hadi Esmaeilzadeh

Alternative Computing Technologies (ACT) Lab  
School of Computer Science, Georgia Institute of Technology

<sup>†</sup>School of Computer Science, Institute for Research in Fundamental Sciences (IPM)

{a.yazdanbakhsh, jspark, hsharma}@gatech.edu      plotfi@ipm.ir      hadi@cc.gatech.edu

## ABSTRACT

General-purpose computing on graphics processing units (GPGPU) accelerates the execution of diverse classes of applications, such as recognition, gaming, data analytics, weather prediction, and multimedia. Many of these applications are amenable to approximate execution. This application characteristic provides an opportunity to improve the performance and efficiency of GPGPU. Recent work has shown significant gains with neural approximate acceleration for CPU workloads. This work studies the effectiveness of neural approximate acceleration for GPU workloads. As applying CPU neural accelerators to GPUs leads to high area overhead, we define a low overhead neurally accelerated architecture for GPGPUs that enables scalable integration of neural acceleration on the large number of GPU cores. We also devise a mechanism that controls the tradeoff between the quality of results and the benefits from neural acceleration. We evaluate this design on a modern GPU architecture using a diverse set of benchmarks. Compared to the baseline GPGPU architecture, the cycle-accurate simulation results show  $2.4\times$  average speedup and  $2.8\times$  average energy reduction with 10% quality loss across all benchmarks. The quality control mechanism retains  $1.9\times$  average speedup and  $2.1\times$  energy reduction while reducing the quality degradation to 2.5%. These benefits are achieved by approximately 1.2% area overhead.

## 1 Introduction

The diminishing benefits from CMOS scaling [1–3] has coincided with the overwhelming increase in rate of data generation. Expert analyses show that in 2011, the amount of generated data surpassed 1.8 trillion GB and by 2020, consumers will generate  $50\times$  this staggering figure [4]. To overcome these challenges, both the semiconductor industry and the research community are exploring new avenues in computer architecture design. Two of the promising approaches are acceleration and approximation. Among programmable accelerators, GPUs provide significant gains in performance and efficiency. GPUs that were originally designed to accelerate graphics functions, now are being used for a wide range of applications, including recognition, learning, gaming, data analytics, weather prediction, molecular dynamics, multimedia, scientific computing, and many more. The availability of programming models for

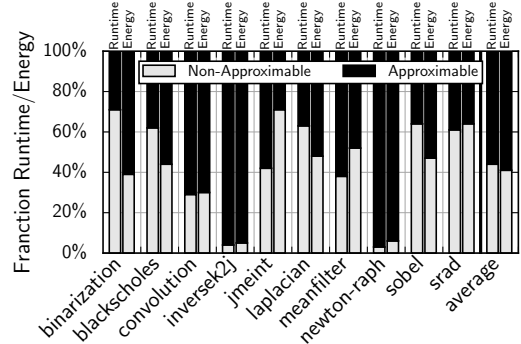


Figure 1: Application runtime and energy breakdown between neurally approximable segments and the regions that cannot be approximated.

general-purpose computing on GPUs and the advances in their microarchitecture has played a significant role in their widespread adoption. Many companies, such as Microsoft, Google, and Amazon use GPUs to provide enterprise services. As GPUs play a major role in executing many classes of applications, improving their performance and efficiency is imperative in enabling new capabilities and coping with the ever-increasing rate of data generation.

Many of the applications that benefit from GPGPUs are also amenable to imprecise computation [6–9]. For these applications, some variation in output is acceptable and some degradation in the output quality is tolerable. This characteristic of many GPU applications provides a unique opportunity to devise approximation techniques that trade small losses in quality for significant gains in performance and efficiency. Neural acceleration is a hardware approximation technique that provides significant gains for CPUs [10–16]. Neural acceleration relies on an automated algorithmic transformation that converts an approximable segment of code<sup>1</sup> to a neural network. This transformation is called the neural transformation [10]. The compiler automatically performs the neural transformation and replaces the approximable segment with an invocation of a neural hardware that accelerates the execution of the thread.

To examine the potential benefits of neural acceleration

<sup>1</sup>Approximable code is a segment that if approximated will not lead to catastrophic failures in execution (e.g., segmentation fault) and its approximation may only lead to graceful degradation of the application output quality.

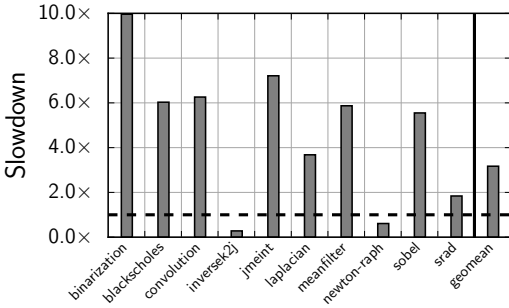


Figure 2: Slowdown with neural transformation due to the lack of hardware support for neural acceleration.

for GPGPU applications, we study<sup>2</sup> its applicability to a diverse set of representative CUDA benchmarks. Figure 1 illustrates the result of this study and shows the breakdown of application runtime and energy dissipation between regions that cannot be neurally approximated and neurally approximable segments<sup>3</sup>. The neurally approximable segments are the ones that can be approximated by a neural network. On average, applications spend 56% of their runtime and 59% of their energy in neurally approximable regions. Some applications such as *inversed2j* and *newton-raph* spend more than 93% of their runtime and energy in neurally approximable regions. These encouraging results demonstrate the significant potential of neural acceleration for GPGPU processors.

**Why not software implementation?** As previous work [17] suggested, it is possible to apply neural transformation with no hardware modifications and replace the approximable segment with an efficient software implementation of the neural network that approximates the region. We explored this possibility and the results are presented in Figure 2. On average, the applications suffer from  $3.2\times$  slowdown. Only *inversed2j* and *newton-raph*, which spent more than 93% of their time in the neurally approximable region, see  $3.6\times$  and  $1.6\times$  speedup, respectively. The slowdown of software implementation is due to (1) the overhead of fetching/decoding the instructions, (2) the overhead of executing the sigmoid function, and (3) the cost of frequent accesses to the memory/register file. The significant potential of neural transformation (Figure 1) and the overall slowdown of software implementation (Figure 2) necessities designing GPU architectures with integrated hardware neural accelerators.

**Why not reusing hardware neural accelerators proposed for CPUs?** Previous work [10] proposes an efficient hardware neural accelerator for CPUs. One possibility is to use CPU neural processing unit (NPU) in GPUs. However, compared to CPUs, GPUs contain (1) significantly larger number of cores (SIMD lanes) that are also (2) simpler. Augmenting each core with an NPU that harbors several parallel processing engines and buffers imposes significant area overhead. Area overhead of integrating NPUs to a GPU while reusing SIMD lanes’ multiply-add units is 31.2%. Moreover, neural networks are structurally parallel. Hence, replacing a code segment with neural networks adds structured parallelism to the thread. In the CPU case, NPU’s multiple multiply-add units exploit this added parallelism to reduce the thread execution latency. GPUs, on the other hand, already exploit data-level parallelism and leverage

<sup>2</sup>Section 6.1 presents our experimental methodology with the GPGPU-Sim cycle-accurate simulator.

<sup>3</sup>The annotation procedure is discussed in Section 2.

many-thread execution to hide thread latencies. One of the insights from this work is that the added parallelism is not the main source of benefits from neural acceleration in GPUs. Therefore, neural acceleration in GPUs leads to a significantly different hardware design as compared to CPUs.

**Contributions.** To this end, the following are the major contributions of this work.

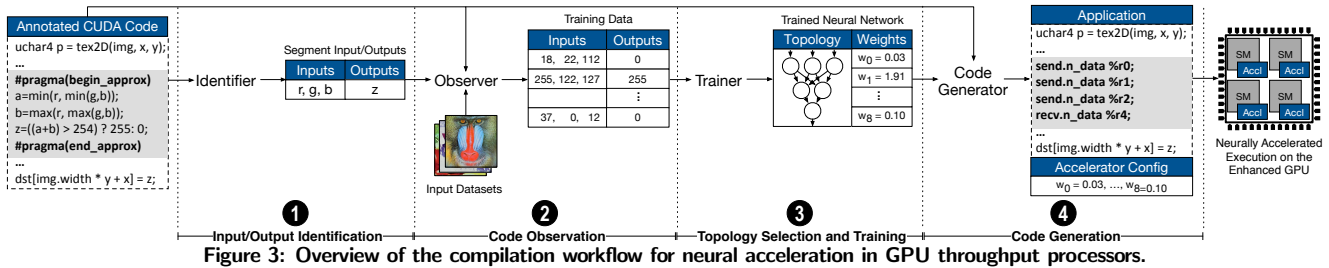
- While this work is not the first to explore neural acceleration, to the best of our knowledge, it is the first to evaluate tight integration of neural acceleration within GPU cores. Integrating neural accelerators within GPUs is fundamentally different than doing so in a CPU because of the hardware constraints and the many-thread SIMT execution model in the GPUs.
- We observe that, unlike in CPUs, the added parallelism is not the main source of benefits from neural acceleration in GPUs. The gains of neural acceleration in GPUs come from (1) storing the parameters and the partial results in small buffers within the SIMD lanes, (2) implementing sigmoid as a lookup table, and (3) eliminating the fetch/decode during neural execution. This insight leads to a low overhead integration of neural accelerators to SIMD lanes by limiting the number of ALUs in an accelerator to only the one that is already in a SIMD lane.
- Through a combination of cycle-accurate simulations and a diverse set of GPU applications from different domains (finance, machine learning, image processing, vision, medical imaging, robotics, 3D gaming, and numerical analysis), we rigorously evaluate the proposed design. Compared to the baseline GPU, our design achieves  $2.4\times$  average speedup and  $2.8\times$  average energy reduction with 10% quality loss. These benefits are achieved with approximately 1.2% area overhead.
- We also devise a mechanism that controls the trade-off between the quality loss and performance and efficiency gains. The quality control mechanism retains  $1.9\times$  average speedup and  $2.1\times$  energy reduction while reducing the quality degradation to 2.5%.

## 2 Neural Transformation for GPUs

To enable integrated neural acceleration on GPUs, the first step is to develop a compilation workflow that can automatically performs the neural algorithmic transformation on GPGPU code. We also need to develop a programming interface that enables developers to delineate approximable regions as candidates for neural transformation. The section describes both the programming interface and the automated compilation workflow for GPGPU applications.

### 2.1 Safe Programming Interface

Any practical approximation technique including ours needs to provide execution safety guarantees. That is, approximation should never lead to catastrophic failures such as out-of-bound memory accesses. In other words, approximation should never affect critical data and operations. The criticality of data and operations is a semantic property of the program and can only be identified by the programmer. The programming language must therefore provide a mechanism for programmers to specify where approximation is safe. This requirement is commensurate with prior work



on safe approximate programming languages such as En-erJ [18], Rely [19], and Axilog [20]. To this end, we extend the CUDA programming language with a pair of `#pragma` annotations that enable marking the start and the end of a safe-to-approximate region of GPGPU applications. The following example illustrates these annotations.

```
#pragma(begin_approx, "min_max")
mi = __min(r, __min(g, b));
ma = __max(r, __max(g, b));
result = ((ma + mi) > 127 * 2) ? 255 : 0;
#pragma(end_approx, "min_max")
```

This segment of the binarization benchmark is approx- imable and is marked as a candidate for transformation. The `#pragma(begin_approx, "min_max")` marks the seg- ment's beginning and names it the "min\_max" segment. The `#pragma(end_approx, "min_max")` marks the end of the segment that was named "min\_max".

## 2.2 Compilation Workflow

As discussed, the main idea of neural algorithmic transfor- mation is to learn the behavior of a code segment using a neural network and then replace the segment with an invo- cation of an efficient neural hardware. To implement this algorithmic transformation, the compiler needs to (1) iden- tify the inputs and outputs of the segment; (2) collect the training data by observing (logging) the inputs and outputs; (3) find and train a neural network that can mimic the ob- served behavior; and finally (4) replace that region of code with instructions that configure and invoke the neural hard- ware. These steps are illustrated in Figure 3. Our compila- tion workflow is similar to the one described in prior work that targets CPU acceleration [10]. However, we specialize these steps for GPGPU applications and add the automatic **input/output identification step to the compilation workflow** to further automate the transformation.

**1 Input/output identification.** To train a neural net- work that mimics a code segment, the compiler needs to collect the input-output pairs that represent the function- ality of the region. The first step is identifying the inputs and outputs of the delineated segment. The compiler uses a combination of live variable analysis and Mod/Ref analy- sis [21] to automatically identify the inputs and outputs of the annotated segment. The inputs are the intersection of live variables at the location of `#pragma(begin_approx)` with the set of variables that are referenced within the segment. The outputs are the intersection of live variables at the location of `#pragma(end_approx)` with the set of variables that are modified within the segment. In the previous example, this analysis identifies `r`, `g`, and `b` as the inputs to the region and `result` as the output.

**2 Code observation.** After identifying the inputs and outputs of the segment, the compiler instruments these

inputs and outputs to log their values in a file as the program runs. The compiler then runs the program with a series of representative input datasets (such as the ones from a program test suite) and logs the pairs of input-output values. The collected set of input-output values constitute the training data that captures the behavior of the segment.

**3 Topology selection and training.** This step needs to both find a topology for the neural network and train it. In finding the topology, the objective is to strike a balance between network's accuracy and its efficiency. Theoretically, a larger, more complex network offers better accuracy potential but is likely to be slower and less efficient. The accuracy of the network does not improve beyond a certain point even if it is enlarged. As follows, the compiler con- sideres a search space for the neural topology and picks the smallest network that delivers comparable accuracy to the largest network in the space. The neural network of choice is Multilayer Perceptron (MLP) that consists of a fully- connected set of neurons organized into layers: the input layer, a number of hidden layers, and the output layer. The number of neurons in the input and output layers is fixed and corresponds to the number of inputs and outputs of the code segment. The problem is finding the number of hidden layers and the number of neurons in each hidden layer.

The space of possible topologies is infinitely large. There- fore, we restrict the search space to neural networks with at most two hidden layers. The number of neurons per hidden layer is also restricted to powers of two, up to 32 neurons. These choices limit the search space to 30 possible topologies. The maximum number of hidden layers and maximum neurons per hidden layer are compilation options and can be changed if needed. These neural networks are trained independently in parallel. To find the best fitting neural network topology, we randomly partition the ap- plication input datasets into a training dataset ( $\frac{2}{3}$  of the programmer-provided application input datasets), and a selection dataset, (the remaining  $\frac{1}{3}$ ). The training datasets are used during training, and the selection datasets are used to select the final neural network topology based on the ap- plication quality loss. Note that we use completely separate input datasets to measure the final quality loss in Section 6.

To train the networks for digital neural acceleration, we use the standard backpropagation [22] algorithm, and for analog neural acceleration, we use the customized learning algorithm presented in [11]. Our compiler performs 10-fold cross-validation for training each neural network. The out- put from this phase consists of a neural network topology – specifying the number of layers and the number of neurons in each layer – along with the weight for each neuron that are determined by the training algorithm.

**4 Code generation.** After identifying the neural network and training it, the compiler replaces the code segment with

special instructions to send the inputs to the neural accelerator and retrieve the results. The compiler also configures the neural accelerator. The configuration includes the weights and the schedule of the operations within the accelerator. This information gets loaded into the integrated neural accelerators when the program loads for execution.

### 3 Instruction Set Architecture Design

To enable neural acceleration, the GPU ISA should provide three instructions: (1) one for sending the inputs to the neural accelerator; (2) one for receiving outputs from the neural accelerator; and finally (3) one for sending the accelerator configuration and the weights. To this end, we extend the PTX ISA with the following three instructions:

1. `send.n_data %r`: This instruction sends the value of register `%r` to the neural accelerator as an input.
2. `recv.n_data %r`: This instruction retrieves a value from the accelerator and writes it to the register `%r`.
3. `send.n_cfg %r`: This instruction sends the value of register `%r` to the accelerator. The instruction also informs the accelerator that the value is for configuration.

We use PTX ISA 4.2 which supports vector instructions that can read or write two or four registers instead of one. We take advantage of this feature and introduce two vector versions for each of our instructions. The `send.n_data.v2 {%r0, %r1}` sends two register values to the accelerator and a single `send.n_data.v4 {%r0, %r1, %r2, %r3}` sends the value of four registers to the neural accelerator. The vector versions for `recv.n_data` and `send.n_cfg` have similar semantics. These vector versions can reduce the number of instructions that need to be fetched and decoded to communicate with the neural accelerator. This reduction lowers the overhead of invoking the accelerator and provides more opportunities for speedup and efficiency gains.

As follows, these instructions will be executed in SIMT mode as other GPU instructions. GPGPU applications typically consist of kernels and GPU threads execute the same kernel code. The neural transformation approximates segments of these kernels. That is, each corresponding thread will contain the aforementioned instructions to communicate with the neural accelerator. Each thread only applies different input data to the same neural network. GPU threads are grouped into cooperative thread arrays (a unit of thread blocks). The threads in different thread blocks are independent and can be executed in any order. The thread block scheduler maps them to GPU processing cores called the streaming multiprocessors (SMs). The SM divides threads of a thread block into smaller groups called warps, typically of size 32 threads. All the threads within a warp execute the same instruction in lock-step. That is, the `send.n_data`, `recv.n_data`, and `send.n_cfg` follow the same SIMT model. That is, executing each of these instructions, conceptually, communicates data with 32 parallel neural accelerators. The GPU-specific challenge is designing a hardware neural accelerator that can be replicated many times within the GPU without imposing extensive hardware overhead. A typical GPU architecture, such as Fermi [23], contains 15 SMs, each with 32 SIMD lanes. That is, to support hardware neural acceleration, 480 neural accelerators need to be integrated. The next section describes our design that scales to such large numbers.

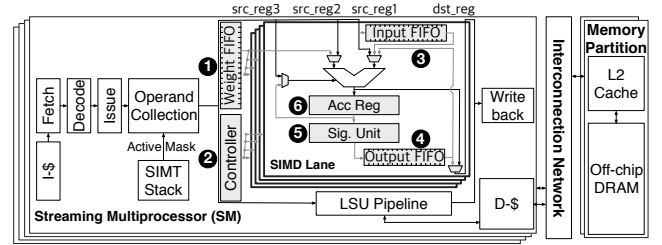


Figure 4: SM pipeline after integrating the digital neural accelerator within SIMD lanes. The added hardware is highlighted by gray.

### 4 Accelerator Design and Integration

To describe our neural accelerator design and its integration into the GPU architecture, we assume a GPU processor based on the Nvidia Fermi. Fermi’s SMs contain 32 double-clock SIMD lanes that execute two half warps (16 threads) simultaneously, where each warp executes in lock-step. Ideally, to preserve the data-level parallelism across the threads and preserve the default SIMT execution model, each SM needs to be augmented with 32 neural accelerators. Therefore, the objective is to design a neural accelerator that can be replicated 32 times with minimal hardware overhead. These two requirements fundamentally change the design space of the neural accelerator from prior work that aims at accelerating single-thread cores with only one accelerator.

A naive approach is to replicate and add a previously proposed CPU neural accelerator to the SMs. These CPU specific accelerators harbor multiple processing engines and contain significant amount of buffering for weights and control. Such a design not only imposes significant hardware overhead, but also is an overkill for data-parallel GPU architectures as our results in Section 6.3 show. Instead, we tightly integrate a GPU specific neural network in every SIMD lane.

Investigating Neural Network Operations As mentioned, the neural algorithmic transformation uses multilayer perceptrons (MLPs) to approximate CUDA code segments. As Figure 5a depicts, an MLP consists of a network of neurons arranged in multiple layers. Each of the neurons in one layer are connected to all of the neurons in the next layer. Each neuron input is associated with a weight value that is the result of training. All neurons are identical and each neuron computes its output ( $y$ ) based on  $y = \text{sigmoid}(\sum_i (w_i \times x_i))$ , where  $x_i$  is a neuron input and  $w_i$  is the input’s associated weight. Therefore, all the computation of a neural network is a set of multiply-add operations followed by the nonlinear sigmoid operation. The neural accelerator only needs to support these two operations.

#### 4.1 Integrating the Neural Accelerator

Each SM has 32 SIMD lanes, divided into two 16-lane groups that execute two half warps simultaneously. The ALU in each lane supports multiply-add. We reuse these ALUs while enhancing the lanes for neural computation. We leverage the SIMT execution to minimize the hardware overhead for the weights and control. We refer to the resulting SIMD lanes as neurally enhanced SIMD lanes.

In Figure 4, the added hardware components are numbered and highlighted in gray. The first component is the **Weight FIFO** (1) that is a circular buffer and stores all of the weights. Since all of the threads are approximated by the same neural network, we only add a **Weight FIFO**, which is shared across all SIMD lanes. The **Weight FIFO** has two read ports corresponding to the two 16 SIMD lanes

that execute two half warps. Each port supplies a weight to 16 ALUs. The second component is the **Controller** (2) which controls the execution of the neural network across SIMD lanes. Again, the **Controller** is shared across 16 SIMD lanes that execute a half warp (two controllers per SM). The **Controller** follows the SIMT pattern of execution for the neural computation and enables the ALUs to perform the computation of the same input of the same neuron in the network.

We augment each of the SIMD lanes with an **Input FIFO** (3) and an **Output FIFO** (4). The **Input FIFO** stores the neural network inputs. The **Output FIFO** stores the output of the neurons including the output neurons that generate the final output. These two are small FIFO structures that are replicated for each SIMD lane. Each of the SIMD lanes also harbors a **Sigmoid Unit** (5) that contains a read-only lookup table, synthesized as combinational logic to reduce the area overhead, that efficiently implements the nonlinear sigmoid function. Finally, the **Acc Reg** (6), which is the accumulator register in each of the SIMD lanes, retains the partial results of the sum of products ( $\sum_i (w_i \times x_i)$ ) before passing it through the **Sigmoid Unit**.

One of the advantages of this design is that it limits all major modifications to SIMD lane pipelines. There is no need to change any other part of the SM except for adding support for decoding the ISA extensions that communicate data to the accelerator (i.e., input and output buffers). Scheduling and issuing these instructions are similar to arithmetic instructions and do not require specific changes.

## 4.2 Executing Neurally Transformed Threads

Figure 5c illustrates the execution of a neurally transformed warp, which contains normal precise and special approximate (i.e., `send.n_data/rcv.n_data`) instructions, on its neurally enhanced SIMD lane pipelines. The other simultaneously executing warp (similarly contains both normal and special instructions) is not shown for clarity. In the first phase (1), SIMD lanes execute the precise instructions as usual before reaching the first `send.n_data` instructions. In the second phase (2), SIMD lanes execute the two `send.n_data` instructions to copy the neural network inputs from the register file to their input buffers. These instructions cause SIMD lanes to switch to the neural mode. In the third phase (3), the enhanced SIMD lanes perform the neural computation and store the results in their output buffers. At the same time, the SM issues `rcv.n_data`, but since the output of the neural network is not ready yet, the SM stops issuing the next instruction and waits for the neurally-enhanced SIMD lanes to finish computing the neural network output. In the fourth phase (4), once the neural network output is ready, `rcv.n_data` instruction copies the results from the output buffer to the register file and then in the fifth phase (5) normal execution resumes. As there is no control divergence or memory access in the neural mode, our design does not swap the running warp with another warp *in the neural mode* to avoid the significant overhead of dedicated input/output buffers or control logic per active warp (SMs support 48 ready-to-execute warps).

## 4.3 Orchestrating Neurally Enhanced Lanes

To efficiently execute neural networks on the neurally enhanced SIMD lanes, the compiler needs to create a static schedule for the neural computation and arrange the weights

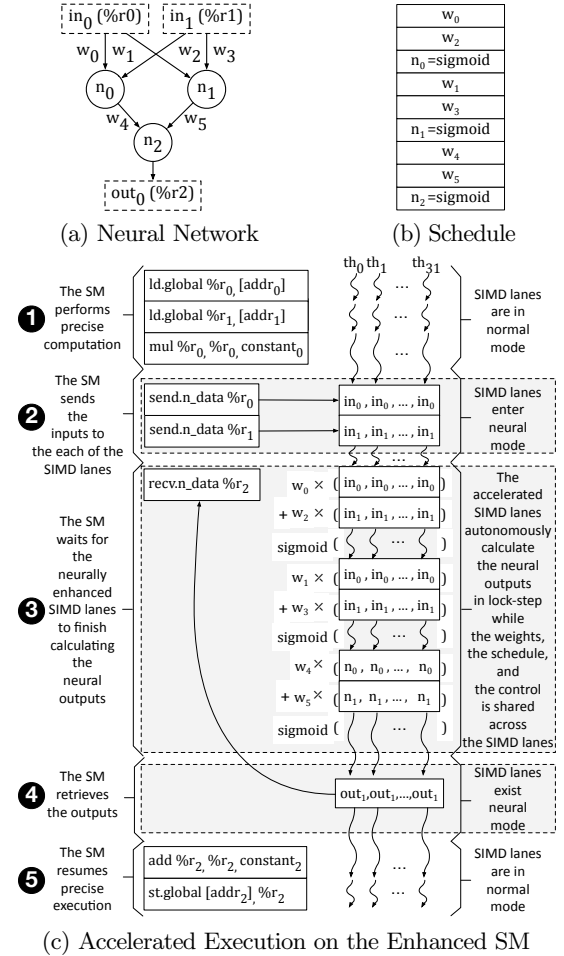


Figure 5: (a) Neural network replacing a segment of a GPU code. (b) Schedule for the accelerated execution of the neural network. (c) Accelerated execution of the GPU code on the enhanced SM.

in proper order. This schedule and the preordered weights are encoded in the program binary and are preloaded to the **Weight FIFO** (Figure 4 1) when the program loads for execution. The compiler generates the execution schedule based on the following steps:

1. The computation of the neurons in each layer has dependence on the output of the neurons in the previous layer. Thus, the compiler first assigns a unique order to the neurons starting from the first hidden layer down to the output layer. This order determines the execution of the neurons. In Figure 5a,  $n_0$ ,  $n_1$ , and  $n_2$  show this order.
2. After ordering the neurons, the compiler generates the order of the multiply-add operations for each neuron, which is followed by a sigmoid operation. This schedule is shown in Figure 5b for the neural network in Figure 5a. The phase 3 of Figure 5c illustrates how the neurally enhanced SIMD lanes execute this schedule in SIMT mode while sharing the weights and control.

The schedule that is presented in Figure 5b constitutes the most of the accelerator configuration and the order in which the weights will be stored in **Weight FIFO** (1 in Figure 4). For each accelerator invocation, SIMD lanes go through these weights in lock-step and perform the neural computation autonomously without engaging the other parts of the SM.

## 5 Controlling Quality Tradeoffs

To be able to control the quality tradeoffs, any approximation technique including ours, needs to expose a quality knob to the compiler and/or runtime system. The knob for our design is the accelerator invocation rate. That is the fraction of the warps that are offloaded to the neural accelerator. The rest of the warps will execute the original precise segment of code and generate exact outputs. In the default case, without any quality control, all the warps that contain the approximable segment will go through the neural accelerator which translates to 100% invocation rate. With quality control, only a fraction of the warps will go through the accelerator. Naturally, the higher the invocation rate, the higher the benefits and the lower the quality.

For a given quality target, the compiler predetermines the invocation rate by examining the output quality loss on a held-out evaluation input dataset. Starts with 100% invocation rate, the compiler gradually reduces the invocation rate until the quality loss is less than the quality target. During the runtime, a quality monitor, similar to the one proposed in SAGE [6], stochastically checks the output quality of the application and adjusts the invocation rate.

We investigated a more sophisticated approach that uses another neural network to filter out invocations of the accelerator that significantly degrade quality. The empirical study suggested that the simpler approach of reducing the invocation rate provides similar benefits.

## 6 Evaluation

We evaluate the benefits of the proposed architecture across different bandwidth and accelerator settings. We use a diverse set of applications, cycle-accurate simulation, logic synthesis, and consistent detailed energy modeling.

### 6.1 Applications and Neural Transformation

**Applications.** As Table 1 shows, we use a diverse set of *approximable* GPGPU applications from the Nvidia SDK [24] and Rodinia [25] benchmark suites to evaluate integrating neural accelerators within GPU architectures. We added three more applications to the mix from different sources [26–28]. As shown, the benchmarks represent workloads from finance, machine learning, image processing, vision, medical imaging, robotics, 3D gaming, and numerical analysis. We did not reject any benchmarks due to their performance, energy, or quality shortcomings.

**Annotations.** As described in Section 2.1, the CUDA source code for each application is annotated using the `#pragma` directives. We use these directives to delineate a region within a CUDA kernel that has fixed number of inputs/outputs and is safe to approximate. Although it is possible and may boost the benefits to annotate multiple regions, we only annotate one region that is easy to identify and is frequently executed. We did not make any algorithmic changes to enable neural acceleration.

As illustrated by the numbers of function calls, conditionals, and loops in Table 1, these regions exhibit a rich and diverse control flow behavior. For instance, the target region in `inversk2j` has three loops and five conditionals. Other regions similarly have several loops/conditionals and function calls. Among these applications, the region in `jmeint` has the most complicated control flow with 37 if/else statements. The regions are also diverse in size and vary

from small (binarization with 27 PTX instructions) to large (`jmeint` with 2,250 PTX instructions).

**Evaluation/training datasets.** As illustrated in Table 1, the datasets that are used for measuring the quality, performance, and energy are completely disjoint from the ones used for training the neural networks. The training inputs are typical representative inputs (such as sample images) that can be found in application test suites. For instance we use the image of `lena`, `peppers`, and `mandrill` for applications that operate on image data. Since the regions are frequently executed, even one application input provides large number of training data. For example, in `sobel` a  $512 \times 512$  pixel image generates 262,144 training data elements.

**Neural networks.** The “Topology” column shows the topology of the neural network that replaces the region of code. For instance, the neural topology for `blackscholes` is  $6 \rightarrow 8 \rightarrow 1$ . That is the neural network has 6 inputs, one hidden layer with 8 neurons, and 1 output neuron. These topologies are automatically discovered by our search algorithm and we use the 10-fold cross validation to train the neural networks. As the results suggest, different applications require different topologies. Therefore, the SM architecture should be changed in a way that is reconfigurable and can accommodate different topologies.

**Quality.** We use an application-specific quality metric, shown in Table 1, to assess the quality of each application’s output after neural acceleration. In all cases, we compare the output of the original precise application to the output of the neurally approximated application. For `blackscholes`, `inversek2j`, `newton-raph`, and `srad` that generate numeric outputs, we measure the average relative error. For `jmeint` that determines whether two 3D triangles intersect, we report the misclassification rate. The convolution, binarization, laplacian, meanfilter, and `sobel` that produce image outputs, we use the average root-mean-square image difference. In Table 1, the “Quality Loss” columns reports the whole-application quality degradation based on the above metrics. This loss includes the accumulated errors due to repeated execution of the approximated region. The quality loss in Table 1 represents the case where all of the dynamic threads with the target region are approximated.

Even with 100% approximation rate, the quality loss with digital neural acceleration is less than 10% except in the case of `jmeint`. The `jmeint` application’s control flow is very complex and the neural network is not able to capture all the corner cases to achieve below 10% quality loss. These results are commensurate with prior work on CPU-based neural acceleration [11, 15]. Prior work on GPU approximation such as SAGE [6] and Paraprox [7] reports similar quality losses in the default setting. EnerJ [18] and Truffle [29] show less than 10% loss for some applications and even 80% loss for others. Green [30] and loop perforation [31] show less than 10% error for some applications and more than 20% for others. Later, we will discuss how to use the invocation rate to control the quality tradeoffs, and achieve even lower quality losses when desired.

To better illustrate the application quality loss, Figure 6 shows the Cumulative Distribution Function (CDF) plot of the final quality loss for each element of the output. Each application output is a collection of elements – an image consists of pixels; a vector consists of scalars; etc. The loss CDF shows the distribution of output quality loss



Table 1: Applications, accelerated regions, training and evaluation datasets, quality metrics, and approximating neural networks.

|                     | Description                            | Source                | Domain             | Quality Metric  | # of Function Calls | # of Loops | # of ifs/elses | # of PTX Insts. | Training Input Set         | Evaluation Input Set          | Digital NPU             |              | Analog NPU              |              |
|---------------------|--|-----------------------|--------------------|-----------------|---------------------|------------|----------------|-----------------|----------------------------|-------------------------------|-------------------------|--------------|-------------------------|--------------|
|                     |  |                       |                    |                 |                     |            |                |                 |                            |                               | Neural Network Topology | Quality Loss | Neural Network Topology | Quality Loss |
| <b>binarization</b> | Image binarization                     | Nvidia SDK            | Image Processing   | Image Diff      | 1                   | 0          | 1              | 27              | Three 512x512 pixel images | Twenty 2048x2048 pixel images | 3 -> 4 -> 2 -> 1        | 8.23%        | 3 -> 8 -> 4 -> 1        | 11.43%       |
| <b>blackscholes</b> | Option pricing                         | Nvidia SDK            | Finance            | Avg. Rel. Error | 2                   | 0          | 0              | 96              | 8,192 options              | 262,144 options               | 6 -> 8 -> 1             | 4.35%        | 6 -> 8 -> 4 -> 1        | 8.23%        |
| <b>convolution</b>  | Data filtering operation               | Nvidia SDK            | Machine Learning   | Avg. Rel. Error | 0                   | 2          | 2              | 886             | 8,192 data points          | 262,144 data points           | 17 -> 2 -> 1            | 5.25%        | 17 -> 4 -> 4 -> 1       | 9.29%        |
| <b>inversek2j</b>   | Inverse kinematics for 2-joint arm     | CUDA-Based Kinematics | Robotics           | Avg. Rel. Error | 0                   | 3          | 5              | 132             | 8,192 2D coordinates       | 262,144 2D coordinates        | 2 -> 16 -> 3            | 8.73%        | 2 -> 16 -> 4 -> 3       | 10.25%       |
| <b>jmeint</b>       | Triangle intersection detection        | jMonkey Game          | 3D Gaming          | Miss Rate       | 4                   | 0          | 37             | 2,250           | 8,192 3D coordinates       | 262,144 3D coordinates        | 18 -> 8 -> 2            | 17.32%       | 18 -> 16 -> 4 -> 2      | 19.70%       |
| <b>laplacian</b>    | Image sharpening filter                | Nvidia SDK            | Image Processing   | Image Diff      | 0                   | 2          | 1              | 51              | Three 512x512 pixel images | Twenty 2048x2048 pixel images | 9 -> 2 -> 1             | 6.01%        | 9 -> 4 -> 2 -> 1        | 9.87%        |
| <b>meanfilter</b>   | Image smoothing filter                 | Nvidia SDK            | Machine Vision     | Image Diff      | 0                   | 2          | 1              | 35              | Three 512x512 pixel images | Twenty 2048x2048 pixel images | 7 -> 4 -> 1             | 7.06%        | 7 -> 8 -> 2 -> 1        | 9.21%        |
| <b>newton-raph</b>  | Newton-Raphson equation solver         | Likelihood Estimators | Numerical Analysis | Avg. Rel. Error | 2                   | 2          | 1              | 44              | 8,192 cubic equations      | 262,144 cubic equations       | 5 -> 2 -> 1             | 3.08%        | 5 -> 4 -> 2 -> 1        | 11.23%       |
| <b>sobel</b>        | Edge detection                         | Nvidia SDK            | Image Processing   | Image Diff      | 0                   | 2          | 1              | 86              | Three 512x512 pixel images | Twenty 2048x2048 pixel images | 9 -> 4 -> 1             | 5.45%        | 9 -> 8 -> 4 -> 1        | 8.03%        |
| <b>srad</b>         | Speckle reducing anisotropic diffusion | Rodinia               | Medical Imaging    | Image Diff      | 0                   | 0          | 0              | 110             | Three 512x512 pixel images | Twenty 2048x2048 pixel images | 5 -> 4 -> 1             | 7.43%        | 5 -> 8 -> 2 -> 1        | 9.87%        |

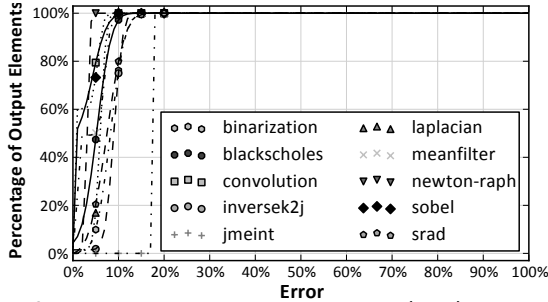


Figure 6: Cumulative distribution function (CDF) plot of the applications output quality loss. A point  $(x, y)$  indicates that  $y$  fraction of the output elements see quality loss less than or equal to  $x$ .

Table 2: GPU microarchitectural parameters.

|  |
|--|
| <b>Processor:</b> 1.4 GHz, No. of SMs: 15, Warp Size: 32 threads/warp, SIMD Width: 8, Max. No. of Threads per Core: 1536, No. of Registers: 32,768, <b>Interconnect:</b> 1 crossbar/direction (15 SMs, 6 MCs), 1.4 GHz <b>L1 Data Cache:</b> 16KB, 128B line, 4-way, LRU; <b>Shared Memory:</b> 48KB, 32 banks; <b>L2 Unified Cache:</b> 768KB, 128B line, 16-way, LRU; <b>Memory:</b> 6 GDDR5 Memory Controllers, 924 MHz, FR-FCFS, <b>Bandwidth:</b> 177.4 GB/sec. |
|--|

among the output elements and shows that very few output elements see a large loss. As shown, the majority of output elements (from 78% to 100%) see a loss less than 10%

## 6.2 Experimental Setup

**Cycle-accurate simulations.** We use the GPGPU-Sim cycle-accurate simulator version 3.2.2 [32]. We modified the simulator to include our ISA extensions and include the extra microarchitectural modifications necessary for integrating neural acceleration within the GPU. The overhead of the extra instructions that communicate the data are modeled in our simulations. For baseline simulations that do not include any approximation or acceleration, we use the unmodified GPGPU-Sim. We use one of GPGPU-Sim’s default configurations that closely models an Nvidia GTX 480 chipset with Fermi architecture. Table 2 summarizes the microarchitectural parameters of the chipset. We also run the applications to completion. We use NVCC 4.2 with -O3 to enable aggressive compiler optimizations. Furthermore, we optimize the number of thread blocks and number of threads per block of each kernel for our simulated hardware.

**Energy modeling and overheads.** To measure the energy benefits, we use GPUWatch [33], which is integrated with GPGPU-Sim. We also generate the event log of the neural accelerator during the cycle-accurate simulations to measure the energy of the neural acceleration. Our energy evaluations use a 40 nm process node and 1.4GHz clock frequency. Neural acceleration requires the following changes to the SM and SIMD lane microarchitecture that are modeled using McPAT [34] and results from CACTI 6.5 [35]. In each SM, we add a 2 KB weight FIFO. The extra input/output FIFO’s are 256 bytes per SIMD lane. The sigmoid LUT which is added to each SIMD lane contains 2048 32-bit entries. Since GPUWatch uses the results from McPAT and CACTI, our added energy models that use the same tools provide a unified and consistent framework for energy measurement.

## 6.3 Experimental Results

**Performance and energy benefits.** Figure 7a shows the whole application speedup when all the invocations of the approximable region are accelerated with the Digital Neural Accelerator (DNA). The remaining part (i.e., the non-approximable region) is executed normally on the GPU. The results are normalized to the baseline where the entire application is executed on the GPU with no acceleration. The highest speedup is observed for *newton-raph* (14.3 $\times$ ) and *inversek2j* (9.8 $\times$ ), where the bulk of execution time is spent on approximable parts (see Figure 1). The lowest speedup is observed for *blackscholes* and *srad* (about 2% and 5%) which are bandwidth-hungry applications.

While a considerable part of the execution time of *blackscholes* and *srad* is spent on approximable parts (See Figure 1), the speedup of accelerating these two applications is minimal because these applications use most of the off-chip bandwidth, even when they run on GPU (without acceleration). Due to bandwidth limitation, DNA accelerators cannot reduce the execution time. Below, we study the effect of increasing the off-chip bandwidth on these two applications and show that with reasonable improvement in bandwidth, even these benchmarks observe significant benefits. On average, the evaluated applications see a 2.4 $\times$

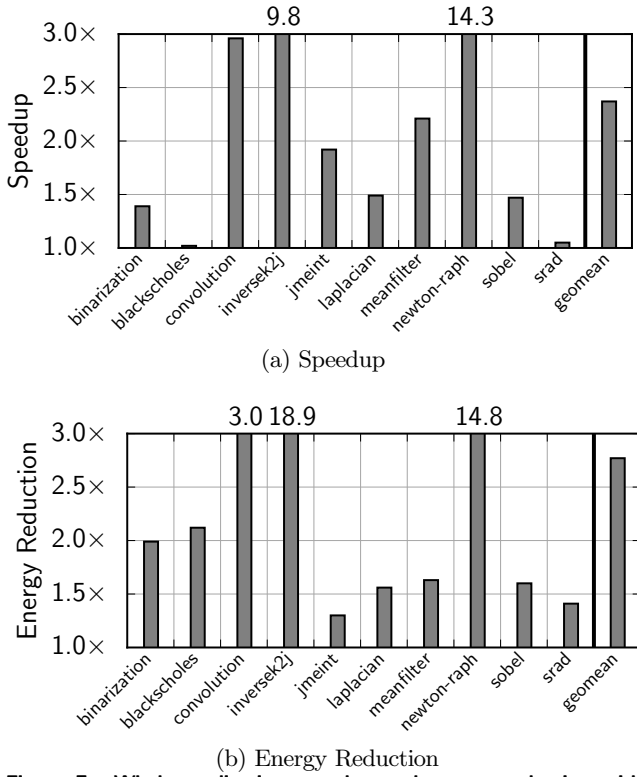


Figure 7: Whole application speedup and energy reduction with GPU+DNA.

speedup through digital neural acceleration.

Figure 7b shows the energy reduction for each benchmark as compared to the baseline where the whole benchmark is executed on GPU. Similar to the speedup, the highest energy saving is achieved for *inversed2j* (18.9 $\times$ ) and *newton-raph* (14.8 $\times$ ), where bulk of the energy is consumed for the execution of approximable parts (see Figure 1). The lowest energy saving is obtained on *jmeint* (30%) as for this application, the fraction of energy consumed on approximable parts is relatively small (See Figure 1). On average, the evaluated applications see a 2.8 $\times$  reduction in energy usage.

The quality loss when all the invocations of the approximable region get executed on DNA (i.e., maximum quality loss) has shown in Table 1 (labeled Quality Loss). We study the effects of our quality control mechanism for trading off performance and energy savings for better quality later in this section.

**Area overhead.** To estimate the area overhead, we synthesize the sigmoid unit using Synopsys Design Compiler and NanGate 45 nm Open Cell library, targeting the same frequency as the SMs. We extract the area of the buffers and FIFOs from CACTI. Overall, the added hardware requires about 0.27 mm<sup>2</sup>. We estimate the area of the SMs by inspecting the die photo of GF100 that implements the Fermi architecture. The area of each SM is about 22 mm<sup>2</sup> and the die area is 529 mm<sup>2</sup> with 15 SMs. The area overhead per SM is approximately 1.2% and the total area overhead is 0.77%. The low area overhead is because our architecture uses the same ALUs that are already available in each of the SIMD lanes, shares the weight buffer across the lanes, and implements the sigmoid unit as read-only lookup table, enabling the synthesis tool to optimize its area. This low area overhead confirms the scalability of our design.

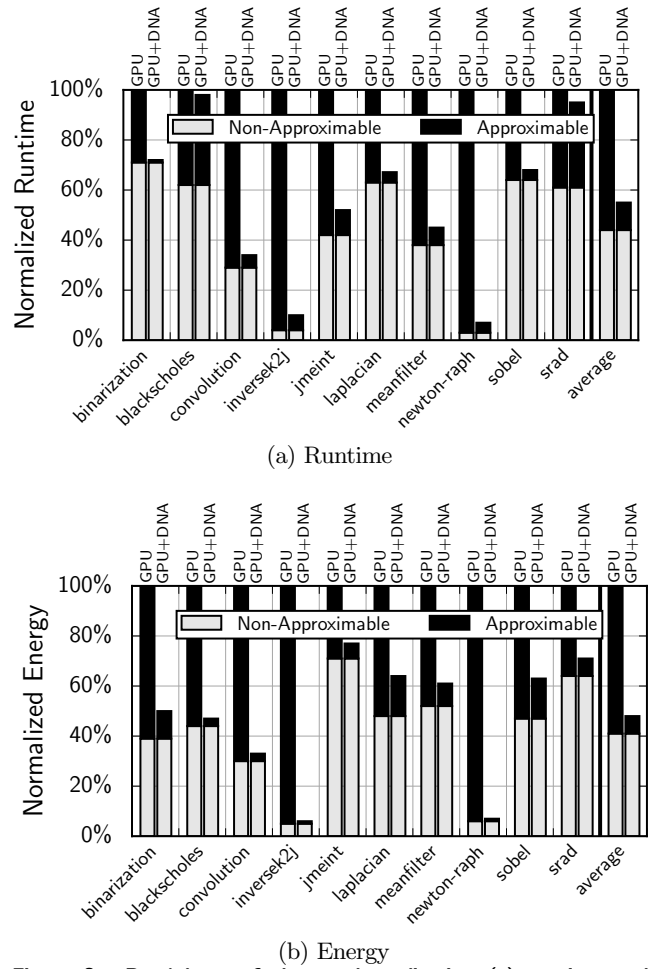


Figure 8: Breakdown of the total application (a) runtime and (b) energy between non-approximable and approximable regions normalized to the runtime and energy consumption of the GPU, respectively. For each application, the first and second bar show the normalized runtime/energy when the application is executed on the GPU, GPU+DNA.

**Opportunity for further improvements.** To explore the opportunity for further improving the execution time by making the neural accelerator faster, Figure 8a shows the time breakdown of approximable and non-approximable parts of applications when applications run on GPU (no acceleration) and GPU+DNA (digital neural acceleration), normalized to the case where the whole application runs on GPU (no acceleration). As Figure 8a depicts, DNA is effective at reducing the time that is spent on approximable parts for all but two applications: *blackscholes* and *sradi*. These two applications use most of the bandwidth of the GPU, and consequently, do not benefit from the accelerators because of hitting the bandwidth wall. The rest of the applications significantly benefit from accelerators.

On some applications (e.g., *binarization*, *laplacian*, and *sobel*), the execution time of approximable parts on DNA is significantly smaller than the execution time of the non-approximable parts. Therefore, there is no further benefit from using a faster accelerator for speeding up the approximable parts. For the rest of the applications, the execution time of approximable parts on DNA (though has reduced considerably) is comparable to (and sometimes exceeds (e.g., *inversed2j*)) the execution time of non-approximable



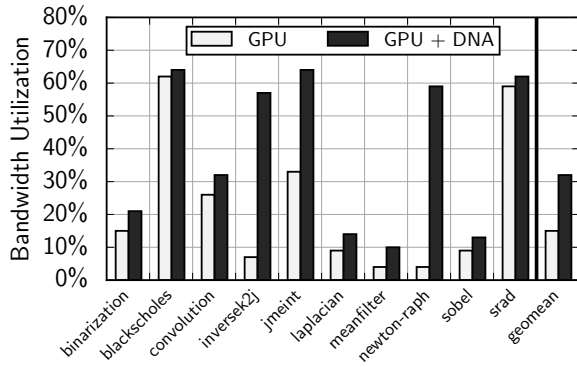


Figure 10: Memory bandwidth consumption when the applications are executed on GPU and GPU+DNA.

parts. Therefore, there is a potential for further speeding up the applications by using faster accelerators.

Likewise, we study the opportunity for further reducing energy usage by benefiting from more energy-efficient accelerators. Figure 8b shows the energy breakdown for approximable and non-approximable parts of applications when applications run on GPU and GPU+DNA, normalized to the case where the whole application runs on GPU. This figure clearly shows that DNA accelerators are extremely efficient at reducing the energy usage of applications on approximable parts. For many of the applications, the energy that is consumed for running approximable parts is insignificant as compared to the energy that is consumed for running the non-approximable parts (e.g., *blackscholes*, *convolution*, *jmeint*, etc.). For these applications, a more energy-efficient neural accelerator implementation brings no further energy saving. However, there are some applications like *binarization*, *laplacian*, and *sobel* for which the fraction of energy that is consumed on DNA accelerators is comparable to the fraction of energy consumed on non-approximable parts. For these applications further energy saving is possible by using a more energy-efficient implementation of neural accelerators.

Below, we first investigate the opportunity of using a faster DNA by varying the speed of the accelerator, and then study the effect of having a more-energy efficient neural network implementation.

**Sensitivity to accelerator speed.** To study the effects of accelerators’ speed on performance gains, we vary the latency of neural accelerators and measure the overall speedup as shown in Figure 9. We decrease the delay of the default accelerators by a factor of 2 and 4 and also include an ideal DNA with zero latency. Moreover, we show the speedup numbers when the latency of the default accelerators increases  $2\times$ ,  $4\times$ ,  $8\times$  and  $16\times$ . Unlike Figure 8a that suggests performance improvement for some applications by benefiting from faster accelerators, Figure 9 shows virtually no speedup benefit by making accelerators faster beyond what they offer in the default design. Even making accelerators slower by a factor of two does not considerably change the speedup. Slowing down the accelerators by a factor of four, many applications observe performance loss. (e.g., *laplacian*).

To explain this behavior, Figure 10 shows the bandwidth usage of GPU and GPU+DNA across all applications. While on the baseline GPU, only two applications use more than 50% of the off-chip bandwidth (i.e., *blackscholes* and *srad*), on GPU+DNA, many applications use more than

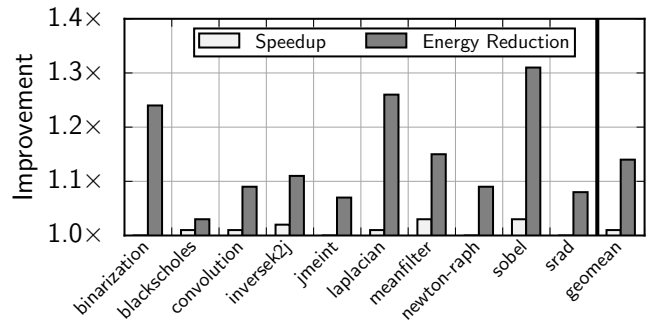


Figure 12: Application speedup and energy reduction with GPU+ANA over GPU+DNA. (The baseline is the accelerated GPU with DNA)

50% of their off-chip bandwidth (e.g., *inversek2j*, *jmeint*, and *newton-raph*). As applications run faster with accelerators, the rate at which they access data increases, which puts pressure on off-chip bandwidth. This phenomena shifts the bottleneck of execution time from computation to data delivery. *As computation is no longer the major bottleneck after acceleration, speeding up thread execution beyond a certain point has marginal effect on the overall execution time.* Even increasing the accelerator speed by a factor of two (e.g., by adding more multiply-and-add units) has marginal effect on execution time. We leverage this insight to simplify the accelerator design and reuse available ALUs in the SMs as described in Section 4.1.

**Sensitivity to off-chip bandwidth.** To study the effects of off-chip bandwidth on the benefits of neural accelerated GPUs, we increase the off-chip bandwidth up to  $8\times$  and report the performance numbers. Figure 11 shows the speedup of GPU+DNA with  $2\times$ ,  $4\times$ , and  $8\times$  bandwidth over the baseline GPU+DNA (i.e.,  $1\times$  bandwidth) across all benchmarks. As GPU+DNA is bandwidth limited for many applications (See Figure 10), we expect a considerable improvement in performance as the off-chip bandwidth increases. Indeed, Figure 11 shows that bandwidth-hungry application (i.e., *blackscholes*, *inversek2j*, *jmeint*, and *srad*) observe speedup of  $1.5\times$  when we double the off-chip bandwidth. After doubling the off-chip bandwidth, no application remains bandwidth limited, and therefore, increasing the off-chip bandwidth to  $4\times$  and  $8\times$  has little effect on performance. It may be possible to achieve, the  $2\times$  extra bandwidth by using data compression [36] with little changes to the architecture of existing GPUs. While technologies like 3D DRAM that offer significantly more bandwidth (and lower access latency) can be useful but are not necessary for providing the off-chip bandwidth requirements of GPU+DNA for the range of applications that we studied. However, even without any of these likely technology advances (compression or stacking), the GPU+DNA provides significant benefits across many applications.

**Analog neural acceleration.** To study the effect of a more energy-efficient implementation of neural accelerators on reducing the energy usage of applications, we evaluate analog implementation of neural accelerators which are more energy efficient than the digital implementation. We use the same design and measurement methodology that is reported in prior work [11]. We use transistor-level SPICE models of the analog neuron. The measurements are from simulation with Cadence Analog Design Environment using predictive technology models at 45 nm [37]. We ran detailed Spectre SPICE simulations. Since the Analog

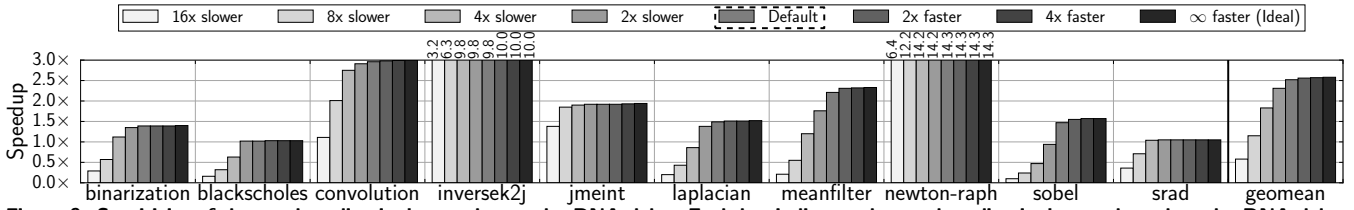


Figure 9: Sensitivity of the total application's speedup to the DNA delay. Each bar indicates the total application's speedup when the DNA delay is altered by different factors. The default delay for DNA varies from one application to the other and depends on the neural network topology trained for that application. The ideal case ( $\infty$  faster) shows the total application speedup when DNA has zero delay.

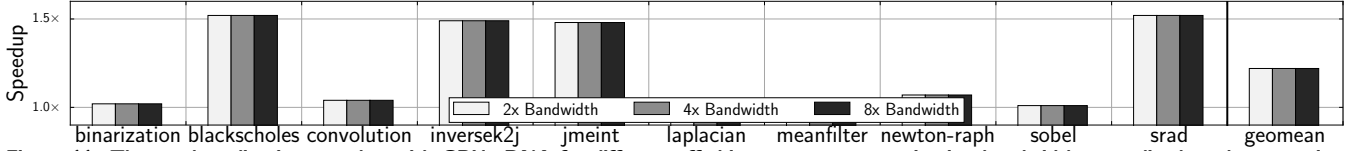


Figure 11: The total application speedup with GPU+DNA for different off-chip memory communication bandwidth normalized to the execution with GPU+DNA with default bandwidth. The default bandwidth is 177.4 GB/s.

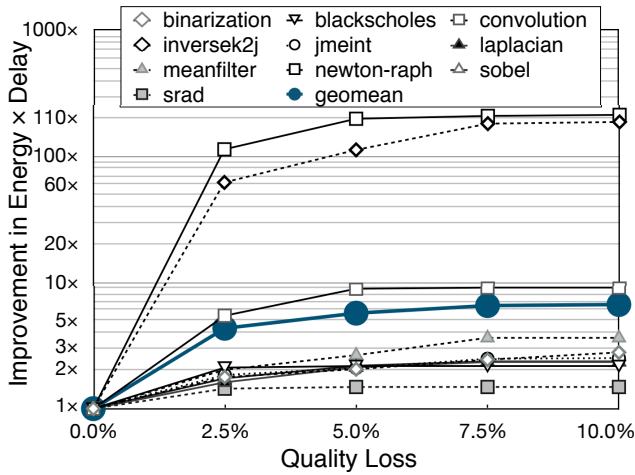


Figure 13: Sensitivity of energy×delay to the output quality.

Neural Accelerators (ANA) cannot use the same ALUs as our digital design, the area overhead is higher. Figure 12 shows the energy and speedup of a GPU with analog neural accelerator (GPU+ANA) as compared to a GPU with digital neural accelerator (GPU+DNA). While ANA is considerably faster than DNA, the speedup of analog neural accelerator matches or slightly exceeds that of the digital implementation across all applications. This is due to the fact that the application speedup does not strictly follow the speed of the accelerators beyond a certain point, as discussed in this section.

However, when it comes to energy saving, some applications benefit from ANA. This is an expected behavior as Figure 8b suggests that some applications benefit from more efficient neural accelerators. While on average the benefit of analog neural accelerators in terms of energy saving is modest, the energy saving on some applications can go as high as  $1.3\times$ . The highest energy saving is observed for binarization, laplacian, and sobel with  $1.2\times$ ,  $1.3\times$  and  $1.3\times$  respectively. These results may not justify the integration of analog acceleration for GPUs. However, it confirms the efficacy of our digital design that can deliver reasonably close benefits to a more energy-efficient analog design.

**Controlling quality tradeoffs.** To study the effect of our quality control mechanism on the gains, Figure 13 shows the energy-delay product of GPU+DNA normalized to

the energy-delay product of the baseline GPU (without acceleration) when the output quality loss changes from 0% to 10%. The proposed quality control mechanism enables navigating the tradeoff between the quality loss and the energy and performance benefits. All of the applications see declines in benefits when invocation rate decreases (i.e., output quality improves). Due to the Amdahl's Law effect, the applications that spend more than 90% of their execution in the approximable segment (inversek2j and newton-raph), see larger declines in benefits when invocation rate decreases. However, even with 2.5% quality loss, the average energy savings is  $2.1\times$  and the average speedup is  $1.9\times$ .

**Comparison with prior CPU neural acceleration.** Prior work [10] has explored improving CPU performance and efficiency with Neural Processing Units (NPUs). Since NPUs offer considerably higher performance and energy efficiency with CPUs, we compare our GPU+DNA proposal to CPU+NPU and GPU+NPU. We use MARSSx86 cycle-accurate simulator for the single-core CPU simulations with a configuration that resembles Intel Nehalem (3.4 GHz with 0.9 V at 45 nm) and is the same as the setup used in the most recent NPU work [11].

Figure 14 shows the application speedup and energy reduction with CPU, GPU, GPU+NPU, and GPU+DNA over CPU+NPU. Even without neural acceleration, GPU provides significant performance and efficiency benefits over NPU-accelerated CPU by leveraging data level parallelism. GPU offers, on average,  $5.6\times$  speedup and  $3.9\times$  energy reduction compared to CPU+NPU. A GPU enhanced with our proposal (GPU+DNA) increases the average speedup and energy reduction to  $13.2\times$  and  $10.8\times$ , respectively. Moreover, as GPUs already exploit data-level parallelism, our proposal offers virtually the same speedup as the area-intensive GPU+NPU. However, accelerating GPU with the NPU design imposes 31.2% area overhead while our GPU+DNA imposes about 1%. GPU with area-intensive NPU offers 17.4% lower energy benefits compared to GPU+DNA mostly due to more leakage. In summary, our proposal offers the highest level of performance and energy efficiency across the examined candidates with the modest area overhead of 1.2% per SM.

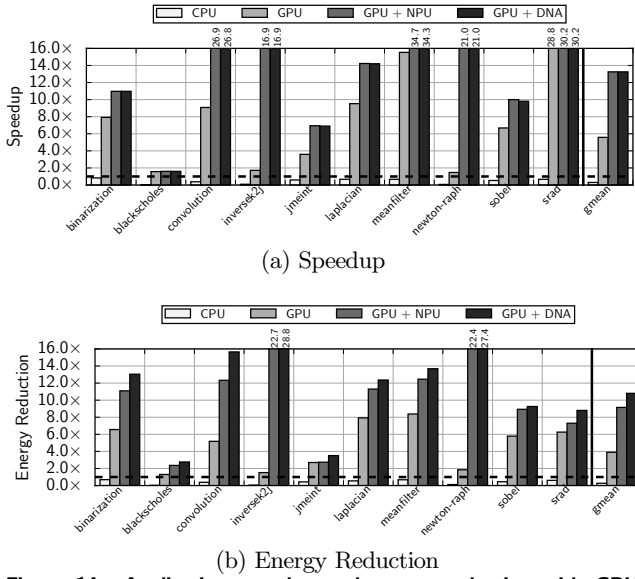


Figure 14: Application speedup and energy reduction with CPU, GPU, GPU+NPU, and GPU+DNA over CPU+NPU [10] (The baseline of CPU+NPU is a CPU with an NPU accelerator).

## 7 Related Work

Recent work has explored a variety of approximation techniques that include: (a) approximate storage designs [38, 39] that trades quality of data for reduced energy [38] and longer lifetime [39], (b) voltage over-scaling [29, 40–43], (c) loop perforation [31, 44, 45], (d) loop early termination [30], (e) computation substitution [6, 9, 30, 46], (f) memoization [7, 8, 47], (g) limited fault recovery [44, 48–53], (h) precision scaling [18, 54], (i) approximate circuit synthesis [20, 55–60], and (j) neural acceleration [10–17].

This work falls in the last category; yet, we exclusively focus on integration of neural accelerators within GPU throughput processors. The prior work on neural acceleration mostly focuses on single-threaded CPU code acceleration by either loosely coupled neural accelerators [12–16] or tightly-coupled ones [10, 11]. Grigorian et al. study the effects of eliminating control-flow divergence by converting SIMD code to software neural networks with no hardware support [17]. However, prior work does not explore tight integration of neural hardware in throughput processors; and does not study the interplay of data parallel execution and hardware neural acceleration. Prior to this work, the benefits, limits, and challenges of integrating hardware neural acceleration within GPUs for many-thread data-parallel applications was unexplored.

There are several other approximation techniques in the literature that can or have been applied to GPU architectures. Loop perforation [31] periodically skips loop iteration for gains in performance and efficiency. Green [30] terminates loops early or substitute compute intensive functions with simpler, lower quality versions that are provided by the programmer. Relax [48] is compiler/architecture system for suppressing hardware fault recovery in approximable regions of code, exposing these errors to the application. Fuzzy memoization forgoes invoking a floating point unit if the inputs are in the neighborhood of previously seen inputs. The results of the previous calculation is reused as an approximate result. Arnau et al. use hardware memoization to reduce redundant computation in GPUs [8].

Sartori et al. propose a technique that mitigates branch divergence by forcing the divergent threads to execute the most popular path [9]. In case of memory divergence, they force all the threads to access the most commonly demanded memory block. SAGE [6] and Praprox [7] perform compile-time static code transformations on GPU kernels that include data compression, profile-directed memoization, thread fusion, and atomic operation optimization. Our quality control mechanism takes inspiration from the quality control in these two works.

In contrast, we describe a hardware approximation technique that integrates neural accelerators within the pipeline of the GPU cores. In our design, we aim at minimizing the pipeline modifications and utilizing existing hardware components. Distinctively, our work explores the interplay between data parallelism and neural acceleration and studies its limits, challenges, and benefits.

## 8 Conclusion

Many of the emerging applications that can benefit from GPU acceleration are amenable to inexact computation. We exploited this opportunity by integrating an approximate form of acceleration, neural acceleration, within the GPU architectures. Our architecture design for the neurally accelerated SMs, provides significant performance and efficiency benefits while providing reasonably low hardware overhead (1.2% area overhead per SM). The quality control knob and mechanism also provided a way to navigate the tradeoff between the quality and the benefits in efficiency and performance. Even with as low as 2.5% loss in quality, our neurally accelerated GPU architecture provides average speedup of 1.9 $\times$  and average energy savings of 2.1 $\times$ . These benefits are more than 10 $\times$  in several cases. These results suggest that *hardware* neural acceleration for GPU throughput processors can be a viable approach to significantly improve their performance and efficiency.

## 9 Acknowledgements

This work was supported by a Qualcomm Innovation Fellowship, NSF award CCF #1553192, Semiconductor Research Corporation contract #2014-EP-2577, and a gift from Google.

## References

- [1] Hadi Esmaeilzadeh, Emily Blem, Renee St. Amant, Karthikeyan Sankaralingam, and Doug Burger. Dark silicon and the end of multicore scaling. In *ISCA*, 2011.
- [2] N. Hardavellas, M. Ferdman, B. Falsafi, and A. Ailamaki. Toward dark silicon in servers. *IEEE Micro*, 31(4):6–15, July–Aug. 2011.
- [3] Ganesh Venkatesh, Jack Sampson, Nathan Goulding, Saturnino Garcia, Vladyslav Bryksin, Jose Lugo-Martinez, Steven Swanson, and Michael Bedford Taylor. Conservation cores: Reducing the energy of mature computations. In *ASPLOS*, 2010.
- [4] John Gantz and David Reinsel. Extracting value from chaos.
- [5] GeForce 400 series. [http://en.wikipedia.org/wiki/GeForce\\_400\\_series](http://en.wikipedia.org/wiki/GeForce_400_series), 2015.

- [6] Mehrzad Samadi, Janghaeng Lee, D Anoushe Jamshidi, Amir Hormati, and Scott Mahlke. SAGE: self-tuning approximation for graphics engines. In *MICRO*, 2013.
- [7] Mehrzad Samadi, Davoud Anoushe Jamshidi, Janghaeng Lee, and Scott Mahlke. Paraprox: pattern-based approximation for data parallel applications. In *ASPLOS*, 2014.
- [8] Jose-Maria Arnau, Joan-Manuel Parcerisa, and Polychronis Xekalakis. Eliminating redundant fragment shader executions on a mobile gpu via hardware memoization. In *Proceeding of the 41st Annual International Symposium on Computer Architecture, ISCA '14*, pages 529–540, Piscataway, NJ, USA, 2014. IEEE Press.
- [9] J. Sartori and R. Kumar. Branch and data herding: Reducing control and memory divergence for error-tolerant gpu applications. *Multimedia, IEEE Transactions on*, 15(2), 2013.
- [10] Hadi Esmaeilzadeh, Adrian Sampson, Luis Ceze, and Doug Burger. Neural acceleration for general-purpose approximate programs. In *MICRO*, 2012.
- [11] Renée St Amant, Amir Yazdanbakhsh, Jongse Park, Bradley Thwaites, Hadi Esmaeilzadeh, Arjang Hassibi, Luis Ceze, and Doug Burger. General-purpose code acceleration with limited-precision analog computation. In *ISCA*, 2014.
- [12] Zidong Du, Avinash Lingamneni, Yunji Chen, Krishna Palem, Olivier Temam, and Chengyong Wu. Leveraging the error resilience of machine-learning applications for designing highly energy efficient accelerators. In *Asia and South Pacific Design Automation Conference (ASP-DAC)*, January 2014.
- [13] Bilel Belhadj, Antoine Joubert, Zheng Li, Rodolphe Hélot, and Olivier Temam. Continuous real-world inputs can open up alternative accelerator designs. In *ISCA*, 2013.
- [14] Beayna Grigorian, Nazanin Farahpour, and Glenn Reinman. BRAINIAC: Bringing reliable accuracy into neurally-implemented approximate computing. In *HPCA*, 2015.
- [15] Thierry Moreau, Mark Wyse, Jacob Nelson, Adrian Sampson, Hadi Esmaeilzadeh, Luis Ceze, and Mark Oskin. SNNAP: Approximate computing on programmable socs via neural acceleration. In *HPCA*, 2015.
- [16] Lawrence McAfee and Kunle Olukotun. Emeuro: A framework for generating multi-purpose accelerators via deep learning. In *Proceedings of the 13th Annual IEEE/ACM International Symposium on Code Generation and Optimization, CGO '15*, pages 125–135, 2015.
- [17] Beayna Grigorian and Glenn Reinman. Accelerating divergent applications on simd architectures using neural networks. In *Computer Design (ICCD), 2014 32nd IEEE International Conference on*, pages 317–323. IEEE, 2014.
- [18] A. Sampson, W. Dietl, E. Fortuna, D. Gnanapragasam, L. Ceze, and D. Grossman. EnerJ: Approximate data types for safe and general low-power computation. In *PLDI*, 2011.
- [19] Michael Carbin, Sasa Misailovic, and Martin C. Rinard. Verifying quantitative reliability for programs that execute on unreliable hardware. In *OOPSLA*, 2013.
- [20] Amir Yazdanbakhsh, Divya Mahajan, Bradley Thwaites, Jongse Park, Anandhavel Nagendrakumar, Sindhuja Sethuraman, Kartik Ramkrishnan, Nishanthi Ravindran, Rudra Jariwala, Abbas Rahimi, Hadi Esmaeilzadeh, and Kia Bazargan. Axilog: Language support for approximate hardware design. March 2015.
- [21] John P Banning. An efficient way to find the side effects of procedure calls and the aliases of variables. In *Proceedings of the 6th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 29–41. ACM, 1979.
- [22] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, volume 1, pages 318–362. MIT Press, 1986.
- [23] Whitepaper: NVIDIA Fermi.
- [24] NVIDIA corporation. NVIDIA CUDA SDK code samples.
- [25] Shuai Che, Michael Boyer, Jiayuan Meng, David Tarjan, Jeremy W Sheaffer, Sang-Ha Lee, and Kevin Skadron. Rodinia: A benchmark suite for heterogeneous computing. In *IISWC*, 2009.
- [26] jMonkeyEngine, 2015.
- [27] Omar Alejandro Aguilar and Joel Carlos Huegel. Inverse kinematics solution for robotic manipulators using a cuda-based parallel genetic algorithm. *Advances in Artificial Intelligence*, 2011.
- [28] Michael Creel and Mohammad Zubair. A high performance implementation of likelihood estimators on gpus. In *CES*, 2013.
- [29] Hadi Esmaeilzadeh, Adrian Sampson, Luis Ceze, and Doug Burger. Architecture support for disciplined approximate programming. In *ASPLOS*, 2012.
- [30] Woongki Baek and Trishul M. Chilimbi. Green: A framework for supporting energy-conscious programming using controlled approximation. In *PLDI*, 2010.
- [31] Stelios Sidiroglou-Douskos, Sasa Misailovic, Henry Hoffmann, and Martin Rinard. Managing performance vs. accuracy trade-offs with loop perforation. In *FSE*, 2011.
- [32] A Bakhoda, G.L. Yuan, W.W.L. Fung, H. Wong, and T.M. Aamodt. Analyzing cuda workloads using a detailed gpu simulator. In *ISPASS*, 2009.

- [33] Jingwen Leng, Tayler Hetherington, Ahmed ElTantawy, Syed Gilani, Nam Sung Kim, Tor M Aamodt, and Vijay Janapa Reddi. Gpuwattch: Enabling energy optimizations in gpgpus. In *Proceedings of the 40th Annual International Symposium on Computer Architecture*, pages 487–498. ACM, 2013.
- [34] Sheng Li, Jung Ho Ahn, Richard D. Strong, Jay B. Brockman, Dean M. Tullsen, and Norman P. Jouppi. McPAT: An integrated power, area, and timing modeling framework for multicore and manycore architectures. In *MICRO*, 2009.
- [35] Naveen Muralimanohar, Rajeev Balasubramanian, and Norm Jouppi. Optimizing NUCA organizations and wiring alternatives for large caches with CACTI 6.0. In *MICRO*, 2007.
- [36] Nandita Vijaykumar, Gennady Pekhimenko, Adwait Jog, Abhishek Bhowmick, Rachata Ausavarungnirun, Onur Mutlu, Chita Das, Mahmut Kandemir, and Todd C. Mowry. A case for core-assisted bottleneck acceleration in gpus: Enabling efficient data compression. In *ISCA*, 2015.
- [37] Yu Cao. Predictive technology models, 2013.
- [38] Song Liu, Karthik Pattabiraman, Thomas Moscibroda, and Benjamin G. Zorn. Flicker: Saving refresh-power in mobile devices through critical data partitioning. In *ASPLOS*, 2011.
- [39] Adrian Sampson, Jacob Nelson, Karin Strauss, and Luis Ceze. Approximate storage in solid-state memories. In *MICRO*, 2013.
- [40] Lakshmi N. Chakrapani, Bilge E. S. Akgul, Suresh Cheemalavagu, Pinar Korkmaz, Krishna V. Palem, and Balasubramanian Seshasayee. Ultra-efficient (embedded) SOC architectures based on probabilistic CMOS (PCMO) technology. In *DATE*, 2006.
- [41] Sriram Narayanan, John Sartori, Rakesh Kumar, and Douglas L. Jones. Scalable stochastic processors. In *DATE*, 2010.
- [42] Rajamohana Hegde and Naresh R. Shanbhag. Energy-efficient signal processing via algorithmic noise-tolerance. In *ISLPED*, 1999.
- [43] Larkhoon Leem, Hyungmin Cho, Jason Bau, Quinn A. Jacobson, and Subhasish Mitra. ERSA: Error resilient system architecture for probabilistic applications. In *DATE*, 2010.
- [44] Sasa Misailovic, Stelios Sidiroglou, Hank Hoffman, and Martin Rinard. Quality of service profiling. In *ICSE*, 2010.
- [45] Martin Rinard, Henry Hoffmann, Sasa Misailovic, and Stelios Sidiroglou. Patterns and statistical analysis for understanding reduced resource computing. In *Onward!*, 2010.
- [46] Jason Ansel, Cy Chan, Yee Lok Wong, Marek Olszewski, Qin Zhao, Alan Edelman, and Saman Amarasinghe. Petabricks: a language and compiler for algorithmic choice. In *PLDI*, 2009.
- [47] Carlos Alvarez, Jesus Corbal, and Mateo Valero. Fuzzy memoization for floating-point multimedia applications. *IEEE Trans. Comput.*, 54(7), 2005.
- [48] Marc de Kruijf, Shuou Nomura, and Karthikeyan Sankaralingam. Relax: An architectural framework for software recovery of hardware faults. In *ISCA*, 2010.
- [49] Xuanhua Li and Donald Yeung. Application-level correctness and its impact on fault tolerance. In *HPCA*, 2007.
- [50] Xuanhua Li and Donald Yeung. Exploiting application-level correctness for low-cost fault tolerance. *J. Instruction-Level Parallelism*, 2008.
- [51] M. de Kruijf and K. Sankaralingam. Exploring the synergy of emerging workloads and silicon reliability trends. In *SELSE*, 2009.
- [52] Yuntan Fang, Huawei Li, and Xiaowei Li. A fault criticality evaluation framework of digital systems for error tolerant video applications. In *ATS*, 2011.
- [53] Vicky Wong and Mark Horowitz. Soft error resilience of probabilistic inference applications. In *SELSE*, 2006.
- [54] Swagath Venkataramani, Vinay K. Chippa, Srimat T. Chakradhar, Kaushik Roy, and Anand Raghunathan. Quality programmable vector processors for approximate computing. In *MICRO*, 2013.
- [55] Ashish Ranjan, Arnab Raha, Swagath Venkataramani, Kaushik Roy, and Anand Raghunathan. Aslan: Synthesis of approximate sequential circuits. In *DATE*, 2014.
- [56] Swagath Venkataramani, Amit Sabne, Vivek Kozhikkottu, Kaushik Roy, and Anand Raghunathan. Salsa: Systematic logic synthesis of approximate circuits. In *DAC*, 2012.
- [57] Jin Miao, A. Gerstlauer, and M. Orshansky. Approximate logic synthesis under general error magnitude and frequency constraints. In *ICCAD*, 2013.
- [58] Kumud Nepal, Yueting Li, R. Iris Bahar, and Sherief Reda. ABACUS: A technique for automated behavioral synthesis of approximate computing circuits. In *DATE*, 2014.
- [59] Avinash Lingamneni, Christian Enz, Krishna Palem, and Christian Piguet. Synthesizing parsimonious inexact circuits through probabilistic design techniques. *ACM Trans. Embed. Comput. Syst.*, 12(2s), 2013.
- [60] Avinash Lingamneni, Kirthi Krishna Muntimadugu, Christian Enz, Richard M. Karp, Krishna V. Palem, and Christian Piguet. Algorithmic methodologies for ultra-efficient inexact architectures for sustaining technology scaling. In *CF*, 2012.